

行動価値関数

(1)ある状態で、特定の行動を選んだ時に、その後どれくらい報酬が期待できるかを示す関数

(2)式:
$$q_{\pi}(s, a) = \sum_{s'} P(s'|s, a)[R(s, a, s') + r v_{\pi}(s')]$$

$q_{\pi}(s, a)$	状態sで行動aを選んだ場合の価値
$\sum_{s'}$	行動aを取った後に遷移する全ての次の状態s'に対して合計
$P(s' s, a)$	遷移確率
r	割引率
$v_{\pi}(s')$	次の状態s'の価値



状態sで行動aを選んだ場合の価値=行動aを取ったあとに遷移するすべての次の状態s'に対しての合計(遷移確率×[即時報酬+割引率×次の状態s'の価値])

行動価値関数とベルマン方程式

(1)ある状態 s で特定の行動 a をとったときに、その後どれくらいの報酬が期待できるかを表す関数

(2)式: $q_{\pi}(s, a) = \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma \sum_{a'} \pi(a'|s')q_{\pi}(s', a')]$

$q_{\pi}(s, a)$	状態 s で行動 a を選んだ場合の価値
$\sum_{s'}$	行動 a をとったあとに遷移するすべての次の状態 s' に対して合計
$P(s' s, a)$	遷移確率
$R(s, a, s')$	即時報酬
γ	割引率
$\sum_{a'}$	次の状態 s' で選べるすべての行動 a' について合計
$\pi(a' s')$	次の状態 s' において、行動 a' を選ぶ確率
$q_{\pi}(s', a')$	次の状態 s' で行動 a' を選んだ場合の価値



状態 s で行動 a を選んだ場合の価値 = 行動 a をとったあとに遷移するすべての次の状態 s' に対して合計(遷移確率 × [即時報酬 + 割引率 × 次の状態 s' で選べるすべての行動 a' についての合計 { 次の状態 s' において行動 a' を選ぶ確率 × 次の状態 s' で行動 a' を選んだ場合の価値 }])